TOWARDS GENERALIZED VIDEO QUALITY ASSESS-MENT: A WEAK-TO-STRONG LEARNING PARADIGM Supplementary Material

Anonymous authors

Paper under double-blind review

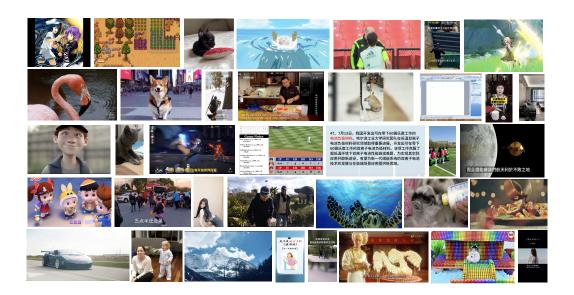


Figure 1: Examples of videos from different categories in our large dataset.

A More Details of Our D_{w2s} Database

A.1 ANALYSIS OF THE COLLECTED VIDEOS

As shown in Fig. 2, our dataset is collected from multiple popular social media platforms with relatively uniform sampling, comprising 20% from Bilibili, 20% from Youku, 25% from YouTube, and 35% from TikTok. All videos are obtained through a filtering pipeline that ensures only publicly available content with permissive licenses is included. Notably, our dataset covers a diverse range of content categories, exceeding twenty in total. In addition to common categories such as lifestyle, food, and animals, it also includes specialized categories such as gaming, Algenerated content, and high-resolution content. To illustrate the diversity of our dataset, we present a variety of video samples in Fig. 1, showcasing the broad range of content available in our large-scale video quality assessment (VQA) dataset. Unlike existing datasets, which often focus on specific formats, our dataset encompasses a wider variety of formats, including both landscape and portrait orientations, as well as various resolutions. This diversity enhances the comprehensiveness of our dataset, making it more suitable for evaluating video quality across a wide kinds of scenarios. A detailed breakdown of our database, including pair types and the corresponding number of videos, is provided in Table 1.

A.2 Analysis of Low-level Metrics

Our data selection strategy is based on a mixed-integer programming method (Vonikakis et al., 2017), which optimizes dataset composition by aligning feature histograms. Specifically, we utilize

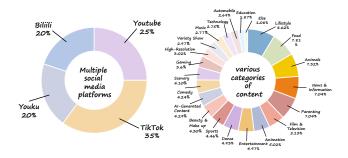


Figure 2: Our dataset is collected from multiple popular social media platforms and encompasses a wide range of content categories.

Table 1: Statistics of raw videos and video pairs in the D_{w2s} dataset.

			_					
Category	Subtype		Videos		Video Pairs			
	V 1	$D_{\mathbf{w2s}}^{(1)}$	$D_{ m w2s}^{(2)}$	$D_{\mathrm{w2s}}^{(3)}$	$D_{\mathrm{w2s}}^{(1)}$	$D_{\rm w2s}^{(2)}$	$D_{{ m w2s}}^{(3)}$	
Ensembling homogeneous teachers	-	200k	100k	50k	250k	85k	85k	
Integrating heterogeneous teachers	Spatial Temporal Compression	50k 20k 10k	2k 1k 1k	2k 1k 1k	160k 40k 50k	5k 5k 5k	5k 5k 5k	
Total		280k	384k	438k	500k	600k	700k	

this approach to match the distributions of nine low-level metrics (blockiness (Romaniak et al., 2012), blur (Narvekar & Karam, 2011), contrast (Peli, 1990), noise, flickering (Pandel, 2008), colourfulness (Hasler & Suesstrunk, 2003), luminance, spatial information (SI) (ITU-T P.910, 2008), and temporal information (TI) (ITU-T P.910, 2008)) between our dataset and the LSVQ dataset. Each metric is computed as follows:

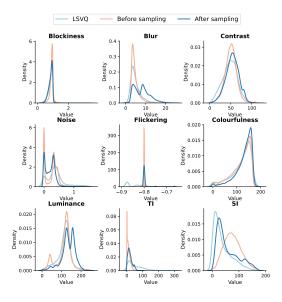


Figure 3: Distribution of nine metrics on the LSVQ dataset, as well as on our dataset before and after sampling.

Blockiness (Romaniak et al., 2012) is quantified by analyzing the luminance differences between pixels within and across encoding blocks. Specifically, we compute the absolute luminance differences between adjacent pixel pairs within the same encoding block (internal pixel pairs) and those spanning adjacent blocks (external pixel pairs). The blockiness metric is then determined as the ratio

 of the total sum of internal pixel difference values to the total sum of external pixel difference values across the entire video frame:

$$B = \frac{\sum_{(x,y)\in\mathcal{I}} |I(x,y) - I(x+1,y)|}{\sum_{(x,y)\in\mathcal{E}} |I(x,y) - I(x+1,y)|},$$
(1)

where I(x,y) represents the luminance value at pixel location (x,y), \mathcal{I} denotes the set of internal pixel pairs, and \mathcal{E} represents the set of external pixel pairs. A higher blockiness value indicates stronger blocking artifacts, which typically result from aggressive video compression.

Blur is measured using the Cumulative Probability of Blur Detection (CPBD) (Narvekar & Karam, 2011), which evaluates perceptual sharpness based on edge width distribution. A higher CPBD value indicates a sharper image. Given an edge pixel e_i , its width $w(e_i)$ is compared with the Just Noticeale Blur (JNB) threshold, determining the blur detection probability $w_{JNB}(e_i)$. The final CPBD score is computed as:

$$CPBD = P(P_{BLUR} \le P_{JNB}) = \sum_{P_{BLUR}=0}^{P_{JNB}} P(P_{BLUR}).$$
 (2)

Contrast is a measure of the dispersion of pixel intensity values within the video frame and can be quantified using the standard deviation of grayscale intensities (Peli, 1990). Specifically, for a grayscale image I(x, y), the mean intensity μ is first computed as:

$$\mu = \frac{1}{M \times N} \sum_{x=1}^{M} \sum_{y=1}^{N} I(x, y), \tag{3}$$

where M and N denote the width and height of the image, respectively, and I(x,y) represents the intensity at pixel (x,y). The contrast value σ is then obtained by calculating the standard deviation of intensity values:

$$\sigma = \sqrt{\frac{1}{M \times N} \sum_{x=1}^{M} \sum_{y=1}^{N} (I(x, y) - \mu)^{2}}.$$
 (4)

The standard deviation σ represents the contrast of the video frame, where a higher σ value indicates a greater dispersion of intensity values and thus a higher contrast.

Flickering occurs when an encoder skips macroblocks to conserve bitrate, especially in low-texture, slow-motion regions (Pandel, 2008). It is quantified by counting macroblock transitions from an "unupdated" to an "updated" state, with a threshold T_f ensuring only significant changes are considered. The flickering metric is computed as:

$$F = \frac{1}{M \times N} \sum_{x=1}^{M} \sum_{y=1}^{N} \mathbb{I}(|I_t(x,y) - I_{t-1}(x,y)| > T_f),$$
 (5)

where $I_t(x,y)$ is the luminance at pixel (x,y) in frame t, and $\mathbb{I}(\cdot)$ is an indicator function. A higher F indicates stronger flickering artifacts.

Colourfulness quantifies color distribution differences across RGB channels, following (Hasler & Suesstrunk, 2003). Given a frame with RGB channels R, G, B, we compute:

$$r_g = R - G, \quad y_b = \frac{1}{2}(R + G) - B.$$
 (6)

The Colourfulness metric is then:

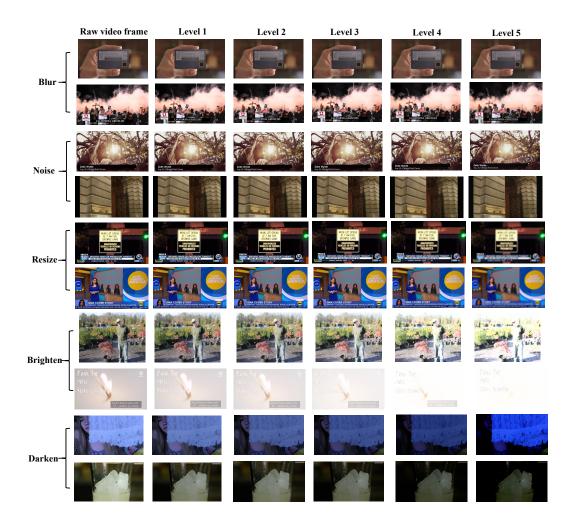


Figure 4: Illustration of different levels of spatial distortion video frames in our large-scale dataset.

$$C = \sqrt{\sigma_{r_g}^2 + \sigma_{y_b}^2} + 0.3 \times \sqrt{\mu_{r_g}^2 + \mu_{y_b}^2},\tag{7}$$

where σ and μ denote the standard deviations and means of r_g and y_b , respectively.

Luminance is measured as the combined intensity of the three RGB channels, defined as:

$$L = R + G + B. (8)$$

SI measures spatial complexity using the Sobel filter. The standard deviation of the Sobel-filtered frame over all pixels is computed, and the maximum value over time represents the SI:

$$SI = \max_{time} \left\{ \operatorname{std}_{space} \left[\operatorname{Sobel}(F_n) \right] \right\}. \tag{9}$$

TI measures motion intensity by calculating the difference between consecutive frames. The temporal difference at pixel (i, j) is:

$$M_n(i,j) = F_n(i,j) - F_{n-1}(i,j). \tag{10}$$

The TI value is the maximum standard deviation of $M_n(i, j)$ over time and space:



Figure 5: Illustration of different levels of streaming distortion video frames in our large-scale dataset.

$$TI = \max_{time} \left\{ \operatorname{std}_{space}[M_n(i,j)] \right\}. \tag{11}$$

To optimize computational efficiency, all metrics are extracted at a sampling rate of one frame per second.

A.3 More Details on Synthetic Distortion Data

A.3.1 SPATIAL DISTORTIONS

We introduce five common spatial distortions: resizing, Gaussian blur, Gaussian noise, darkening, and brightening. Each distortion is applied at five different levels to simulate varying degrees of degradation, ranging from mild to severe. Fig. 4 illustrates examples of these distortions, where the quality of video frames progressively deteriorates as the distortion level increases. Below, we provide details on how these spatial distortions are generated, where I represents the original frame, and I' denotes the distorted frame.

Resizing: The frame is first downsampled by a scaling factor s and then upsampled back to its original size. This process reduces spatial details and introduces pixelation artifacts, simulating resolution loss. The transformation is defined as:

$$I' = \text{Upsample}(\text{Downsample}(I, s), s), \tag{12}$$

where s takes values from the set $\{2, 3, 4, 8, 16\}$.

Gaussian Blur: The frame is convolved with a Gaussian kernel, where the standard deviation σ_{blur} controls the extent of the blur. A larger σ_{blur} results in a wider spread of the Gaussian function, leading to a stronger blurring effect by averaging pixel intensities over a larger neighborhood. The blurring process is defined as:

$$I' = I * G(\sigma_{blur}), \tag{13}$$

where $G(\sigma_{blur})$ is a Gaussian kernel with standard deviation σ_{blur} which takes values from the set $\{0.1, 0.5, 1, 2, 5\}$, and * denotes the convolution operation.

Table 2: An overview of our testing datasets.

Dataset	Year	# of Videos	# of Scenes	Resolution	Duration	Frame Rate	Distortion Type
KoNViD-1k (Hosu et al., 2017)	2017	1,200	1,200	540p	8	24, 25, 30	In-the-wild
LIVE-VQC (Sinno & Bovik, 2018)	2018	585	585	240p-1080p	10	30	In-the-wild
YouTube-UGC (Wang et al., 2019)	2019	1,380	1,380	360p-4K	20	30	In-the-wild
LSVQ (Ying et al., 2021)	2021	38,811	38,811	99p-4K	5-12	< 60	In-the-wild
Waterloo-IVC-4K (Li et al., 2019)	2019	1200	20	540p, 1080p, 4k	9-10	24, 25, 30	H.264 compression
LIVE-YT-HFR (Madhusudana et al., 2021)	2021	480	16	1080p	6-10	24, 30, 60, 82, 98, 120	Frame rate, VP9 compression
LIVE-YT-Gaming (Yu et al., 2022)	2022	600	600	360p-1080p	8-9	30, 60	PGC, UGC
CGVDS (Saha et al., 2023)	2023	360	15	480p, 720p, 1080p	30	20, 30, 60	H.264 compression
KVQ (Lu et al., 2024)	2024	4200	600	-	3-8	-	UGC

Gaussian noise: Gaussian noise is introduced by adding random variations to each pixel, following a normal distribution with mean μ and standard deviation σ_{noise} . The noise level is controlled by adjusting σ_{noise} , where higher values result in more pronounced noise artifacts. The process is defined as:

$$I' = I + N(\mu, \sigma_{noise}^2), \tag{14}$$

where $N(\mu, \sigma_{noise}^2)$ represents Gaussian noise with mean μ and variance σ_{noise}^2 , added independently to each pixel. σ takes values from the set $\{0.001, 0.002, 0.003, 0.005, 0.01\}$.

Darkening: Darkening is applied by reducing the luminance component in the color space. The effect is controlled by a parameter p, which determines the degree of brightness reduction. The luminance channel L is adjusted using an interpolation function f(L, p) as follows:

$$L' = f(L, p). (15)$$

The parameter p is selected from a predefined set of values $\{0.05, 0.1, 0.2, 0.4, 0.8\}$, with larger values leading to stronger darkening effects.

Brightening: In contrast, brightening is achieved by enhancing the luminance component in the color space. The luminance channel L is modified using a nonlinear transformation function g(L, p):

$$L' = g(L, p), \tag{16}$$

The parameter p is selected from $\{0.1, 0.2, 0.4, 0.7, 1.1\}$, with larger values producing a stronger brightening effects.

A.3.2 TEMPORAL DISTORTIONS

We introduce two types of temporal distortions: jitter and stuttering, each distortion maintain three different levels.

Jitter: Jitter introduces random shifts and random cropping followed by resizing of video frames. The amount of shift is determined by the jitter level, which controls the extent of spatial displacement.

For each frame, random horizontal and vertical shifts are applied using an affine transformation matrix, which shifts the frame along the x- and y-axes. Additionally, each frame is cropped by a small amount from the edges and resized back to its original dimensions, simulating pixelation effects or lower-quality views. The transformation matrix is described as follows:

$$M = \begin{bmatrix} 1 & 0 & \text{random_shift_x} \\ 0 & 1 & \text{random_shift_y} \end{bmatrix}$$
 (17)

where random_shift_x and random_shift_y are random values determined by the jitter level.

Stuttering: Stuttering is introduced by randomly dropping frames at a controlled rate. The drop rate p_d is determined by the distortion level, where higher levels correspond to increased frame loss. For each frame I_t , a random probability is drawn and compared with p_d . If the frame is dropped, it is replaced by the previous frame I_{t-1} , simulating temporal freezing in the video. The process can be formulated as:

$$I_t' = \begin{cases} I_{t-1}, & \text{if } r < p_d, \\ I_t, & \text{otherwise} \end{cases}$$
 (18)

where $r \sim U(0,1)$ is a random variable drawn from a uniform distribution.

A.3.3 STREAMING DISTORTIONS

As illustrated in Fig. 5, we select the two most common compression standards, H.264 and H.265, to simulate video quality degradation for the compression distortion. These distortions are applied using the ffmpeg tool, a widely used multimedia framework, to encode the videos with different compression settings. Specifically, we chose four fixed constant rate factor (CRF) values for each compression standard to control the level of distortion.

For H.264 compression, we selected the fast encoding mode, which provides a good balance between encoding speed and compression efficiency, making it suitable for real-time applications. To cover a wide range of compression levels, we applied H.264 compression using CRF values of 24, 36, 48, and 63, ensuring the simulation of various quality degradation scenarios.

In contrast, for H.265 compression, we selected the very slow encoding mode, which prioritizes compression efficiency over speed, leading to higher quality video at the cost of longer encoding times. To achieve fine-grained quality simulation, we applied H.265 compression with a narrower CRF range of 36, 40, 44, and 48, allowing for precise control over compression artifacts.

These encoding settings help to simulate typical real-world compression scenarios, where different modes and CRF values are chosen based on the trade-off between video quality and encoding performance.

A.4 More Details on Testing Datasets

Table 2 provides an overview of our testing datasets, which encompass diverse content types, resolutions, durations, frame rates, and distortion types. The first four datasets consist of in-the-wild videos containing various authentic distortions, while the remaining datasets focus on specific content types and distortion factors. For example, LIVE-YT-Gaming is dedicated to gaming content, LIVE-YT-HFR targets frame rate distortions, and Waterloo-IVC-4K covers different types of compression artifacts. By evaluating our model across these nine datasets, we demonstrate its robustness and effectiveness in both in-domain and out-of-distribution (OOD) quality assessment scenarios.

B MORE DETAILS OF QUALITY ANNOTATION

B.1 WEAK MODELS FOR PSEUDO-LABELING

Table 3: Comparison of model parameters and architecture.

Model	Parameters (M)	Architecture				
MinimalisticVQA(VII)	86.93	Swin-B				
MinimalisticVQA (IX)	121.59	Swin-B + SlowFast				
FAST-VQA	29.97	Swin-Tiny				
DOVER	58.06	Swin-Tiny + Conv-Tiny				
Q-Align	8204.56	mPLUG-Owl2				
Our strong model	8075.24	LLaVA-OneVision-Chat + SlowFas				

We choose five SOTA VQA models: Minimalistic VQA (VII) (Sun et al., 2024), Minimalistic VQA (IX) (Sun et al., 2024), FAST-VQA (Wu et al., 2022), DOVER (Wu et al., 2023a), and Q-Align (Wu et al., 2023b) as weak teachers to formulate our pseudo quality annotation. The detail introduction of the five models is as follows:

MinimalisticVQA (VII) employs Swin Transformer-B (Liu et al., 2022), pre-trained on ImageNet-1K (Deng et al., 2009), as the spatial quality analyzer to extract quality-aware spatial features from key frames, ensuring robust spatial quality assessment.

MinimalisticVQA (**IX**) builds upon MinimalisticVQA (VII) by incorporating a temporal quality analyzer to account for motion distortions. The temporal quality analyzer, implemented using the SlowFast (Feichtenhofer et al., 2019) network pre-trained on the Kinetics-400 (Carreira & Zisserman, 2017) dataset, extracts motion-related features from video chunks, enhancing the model's ability to assess temporal quality variations.

FAST-VQA introduces Grid Mini-patch Sampling (GMS) strategy, which preserves local quality by sampling patches at raw resolution and maintains global quality through uniformly sampled minipatches. These mini-patches are spliced and temporally aligned into fragments. To process these fragments, the Fragment Attention Network (FANet) is designed to effectively extract video quality features. Combining GMS and FANet, FAST-VQA achieves efficient end-to-end video quality assessment with effective feature representation learning.

DOVER builds upon FAST-VQA as its technical branch to capture low-level distortions, while introducing an additional aesthetic branch to assess high-level semantic composition, which relates to user preferences and content recommendation. By disentangling these two perspectives, DOVER establishes a more human-aligned and interpretable framework for video quality assessment.

Q-Align presents a novel training strategy for large multimodal model (LMM) in VQA by replacing direct numerical score predictions with discrete, text-defined rating levels (e.g., "excellent", "good", "fair", "poor", "bad") as learning targets. During inference, Q-Align extracts the log probabilities of each rating level, applies softmax normalization to obtain a probability distribution, and computes a weighted average to derive the final predicted quality score.

B.2 PROMPTS FOR MODEL TRAINING

We construct the label prompts for our large-scale dataset using a fixed template. For the single-video input:

```
Question: "You will now receive a video: <image>. Please watch the video carefully and answer the following question: What is your overall rating of the quality of this video?" Answer: "[quality score]"
```

For the dual-video input:

```
Question: "You will now receive two videos. The first video: <image>. The second video: <image>. Please watch both videos carefully and answer the following question: Compared to the first video, how would you rate the quality of the second video?"

Answer: "The quality of the second video is [level] compared to the first video."
```

C More Details of Our Strong student Model

C.1 MODEL STRUCTURE

As illustrated in Fig. 3 in the main paper, our model comprises three components: a visual feature extractor, a text tokenizer, and an LLM decoder.

Visual Feature Extractor. The visual feature extractor adopts a dual-branch design: a spatial branch with image encoder \mathcal{F}_I (i.e., SigLIP) processes key frames, while a temporal branch with pre-trained

motion encoder \mathcal{F}_M (i.e., SlowFast) analyzes frame sequences. Both branches employ dedicated projection layers $\mathcal{P}_{\mathcal{I}}$ and $\mathcal{P}_{\mathcal{F}}$ (i.e., two-layer MLPs) to map spatial and temporal features into visual tokens aligned with language space. Specifically, given an input video $\boldsymbol{x} = \{\boldsymbol{x}_i\}_{i=0}^{N-1}$ containing N frames at frame rate r, we first partition it into $N_c = \lfloor N/r \rfloor$ continuous chunks $\{\boldsymbol{c}_k\}_{k=0}^{N_c-1}$, where each chunk $\boldsymbol{c}_k = \{x_j\}_{j=k*r}^{(k+1)*r}$ spans r frames. Spatial features \boldsymbol{f}_k^s are extracted from the first frame \boldsymbol{x}_{kr} of each chunk, while temporal features \boldsymbol{f}_k^t are computed over all frames in c_k . The feature extraction process is formally expressed as:

$$f_k^s = \mathcal{P}_I(\mathcal{F}_I(\boldsymbol{x}_{kr})), \quad f_k^t = \mathcal{P}_M(\mathcal{F}_M(\boldsymbol{c}_k)),$$

$$f^v = \operatorname{Concat}\left([\boldsymbol{f}_k^s, \boldsymbol{f}_k^t]_{k=0}^{N_c - 1}\right),$$
(19)

where f^v is the extracted visual features of x. Given a video pair (x^A, x^B) , we can derive the visual features (f_A^v, f_B^v) .

Feature Fusion via the LLM. Given an input prompt p, we first encode it into text tokens $f^p = \mathcal{T}(p)$ using tokenizer \mathcal{T} . The visual features of a video pair (f_A^v, f_B^v) are then concatenated with f^t and fed to a pretrained LLM decoder (i.e., Qwen-2) for multimodal fusion to derive the output response for quality ranking:

$$r = \mathcal{L}(f_A^v, f_B^v, f^p), \tag{20}$$

where r is expected to belong to {"superior", "better", "similar", "worse", "inferior"}.

C.2 TRAINING DETAILS

C.2.1 TRAINING SETUP

The model is trained using the DeepSpeed framework with mixed-precision floating-point operations to optimize memory and computational efficiency. The training is conducted for one epoch with a batch size of 1 per device and a gradient accumulation step of 1. The optimizer follows AdamW with a initial learning rate of 1×10^{-4} , a cosine learning rate schedule, and a warm-up ratio of 0.03.

We employ a joint training strategy for images and videos. For the image encoder, videos are sampled at a rate of one frame per second, with each sampled frame resized to a resolution of 384×384 , while images are directly resized to the same resolution. For the motion encoder, videos are fully encoded across all frames to capture temporal dynamics, whereas images, which lack temporal information, are assigned an all-zero tensor as their temporal representation.

C.2.2 AUXILIARY CONFIDENCE LOSS

As mentioned in Section 4.3 in the main paper, we introduce an auxiliary confidence loss to encourage the model to maintain high-confidence predictions, especially in the presence of noisy weak supervision. The final training objective is a dynamically weighted combination of the cross-entropy loss \mathcal{L}_{CE} and the confidence loss \mathcal{L}_{conf} :

$$\mathcal{L} = (1 - \lambda) \cdot \mathcal{L}_{CE} + \lambda \cdot \mathcal{L}_{conf}, \tag{21}$$

where λ is an adaptive weighting factor that balances between trusting the weak labels and relying on the model's own confidence. The confidence loss is defined as the average entropy over the predicted token probability distributions:

$$\mathcal{L}_{\text{conf}} = \frac{1}{N} \sum_{i=1}^{N} H(p_{\theta}(x_i)) = -\frac{1}{N} \sum_{i=1}^{N} \sum_{c} p_{\theta}(c|x_i) \log p_{\theta}(c|x_i), \tag{22}$$

where $p_{\theta}(c|x_i)$ denotes the predicted probability of vocabulary token c given input x_i . By minimizing the entropy of the predicted distribution, we encourage the model to produce more confident next-token predictions.

To dynamically adjust λ during training, we introduce a temperature-based confidence estimation mechanism. Specifically, we define:

$$\lambda = \alpha \cdot \min\left(1.0, \frac{t}{T_{\text{warmup}}}\right),\tag{23}$$

where t denotes the current training step ratio (normalized to [0,1]), and $T_{\rm warmup}$ is the warm-up period, which we set to 10% of the total training steps. This warm-up phase ensures that the strong model gradually learns to rely on its own confidence, while initially being guided by the weak labels. The factor α is computed as the ratio between the temperature-scaled exponentials of the two losses:

$$\alpha = \frac{\exp(\mathcal{L}_{\text{conf}}/T)}{\exp(\mathcal{L}_{\text{conf}}/T) + \exp(\mathcal{L}_{\text{CE}}/T)}.$$
 (24)

Here, T is a temperature parameter that controls the sharpness of the weighting between the two loss components. We linearly decrease T from 0.5 to 0.1 during the warm-up period to gradually increase the sensitivity of α to differences in the two loss values.

C.3 INFERRING DETAILS

C.3.1 PROBABILITY MODELING

Though we employ video pairs to train our model by enabling it to determine whether the second video is better than the first, our goal during inference is to obtain an absolute quality score for a single video. To achieve this, we propose a method that converts the probability of a test video being better or worse than anchor videos into a final quality score.

First, we describe how to construct the probability distribution for comparative quality assessments. The comparative token set is defined as:

$$S = \{s_k\}_{k=1}^5 = \{inferior, worse, similar, better, superior\}.$$
 (25)

The probability of each token is computed using the softmax function:

$$q_{s_k} = \frac{e^{s_k}}{\sum_{m=1}^r e^{s_m}},\tag{26}$$

where q_{s_k} represents the probability of the k-th token, and r denotes the number of levels.

To obtain a quality score for the test video $v_{\rm eval}$, we aggregate its comparative probabilities against anchor videos using a weighted summation:

$$P(v_{\text{anchor}}, v_{\text{eval}}) = \sum_{k=1}^{r} \alpha_k q_{s_k} (v_{\text{anchor}}, v_{\text{eval}}), \quad r = 1 \dots p.$$
 (27)

where α_k are fixed weights that reflect the comparative levels. Specifically, the weights are defined as:

$$\{\alpha_k\}_{k=1}^5 = \{0, 0.25, 0.5, 0.75, 1\}.$$
 (28)

This approach enables the model to generate a continuous quality score for a single video by leveraging its relative comparisons against anchor videos in the training set.

C.3.2 SCORE MODELING

Finally, we construct a probability matrix based on pairwise comparisons with a set of anchor videos. Given a set of five anchor videos, we first define a probability matrix:

$$M_r \in \mathbb{R}^{5 \times 5},$$
 (29)

where each entry $P(b^{(i)}, b^{(j)})$ represents the probability that anchor video $b^{(i)}$ is preferred over $b^{(j)}$. This probability satisfies:

LSVQ-labeled

Table 4: Performance of weak-to-strong models trained with pseudo-labels from weak models. For comparison, we also report the performance of our model trained directly on the LSVQ dataset.

MinimalisticVQA(IX) 0.885 0.882 0.792 0.828 0.862 0.859 0.775 0.821 0.826 0.821 0.849 0.88 FAST-VQA 0.880 0.880 0.881 0.813 0.859 0.854 0.826 0.845 0.730 0.747 0.838 0.8 Q-Align 0.886 0.884 0.761 0.822 0.876 0.878 0.783 0.819 0.844 0.771 0.771 0.781 0.842 0.8 Weak-to-Strong Students 0.866 0.885 0.762 0.795 0.859 0.857 0.771 0.813 0.808 0.821 0.824 0.8 Minimalistic VQA(IX)-labeled 0.879 0.878 0.794 0.826 0.869 0.871 0.868 0.785 0.819 0.843 0.825 0.824 0.8 DOVER-labeled 0.871 0.869 0.780 0.813 0.870 0.875 0.819 0.813 0.840 0.8 Supervised Student <t< th=""><th>In-domain Datasets</th><th colspan="2" rowspan="2">- Crest</th><th colspan="2" rowspan="2">3,573</th><th colspan="2" rowspan="2">1,200</th><th colspan="2" rowspan="2">LIVE-VQC 585</th><th colspan="2" rowspan="2">YouTube-UGC 1,020</th><th colspan="2" rowspan="2">Overall -</th></t<>	In-domain Datasets	- Crest		3,573		1,200		LIVE-VQC 585		YouTube-UGC 1,020		Overall -	
MinimalisticVQA(VII) 0.861 0.859 0.740 0.784 0.843 0.841 0.757 0.813 0.775 0.779 0.817 0.881 0.885 0.882 0.792 0.828 0.862 0.859 0.775 0.821 0.826 0.821 0.849 0.845 0.845 0.845 0.826 0.821 0.849 0.845 0.845 0.826 0.845 0.730 0.747 0.838 0.880 0.880 0.881 0.813 0.857 0.854 0.826 0.845 0.730 0.747 0.838 0.886 0.884 0.761 0.822 0.876 0.878 0.869 0.817 0.840 0.771 0.781 0.842 0.845 0.826 0.845 0.730 0.747 0.838 0.845 0.826 0.845 0.730 0.747 0.838 0.846 0.845 0.826 0.845 0.826 0.845 0.826 0.845 0.771 0.813 0.842 0.846 0.844 0.845 0.846 0.846 0.8	# of videos												
MinimalisticVQA(VII) 0.861 0.859 0.740 0.784 0.843 0.841 0.757 0.813 0.775 0.779 0.817 0.88 MinimalisticVQA(IX) 0.885 0.882 0.792 0.828 0.860 0.878 0.860 0.781 0.813 0.859 0.854 0.826 0.821 0.820 0.838 0.86 DOVER 0.880 0.880 0.880 0.881 0.813 0.859 0.851 0.845 0.826 0.843 0.747 0.838 0.86 Q-Align 0.886 0.884 0.761 0.822 0.876 0.817 0.841 0.771 0.711 0.812 0.844 0.84 Week-to-Strong Students MinimalisticVQA(IX)-labeled 0.855 0.852 0.762 0.795 0.859 0.857 0.771 0.813 0.804 0.824 0.88 MinimalisticVQA(IX)-labeled 0.879 0.878 0.795 0.819 0.849 0.855 0.789 0.834 0.875 <t< th=""><th>Methods</th><th>SRCC</th><th>PLCC</th><th>SRCC</th><th>PLCC</th><th>SRCC</th><th>PLCC</th><th>SRCC</th><th>PLCC</th><th>SRCC</th><th>PLCC</th><th>SRCC</th><th>PLCC</th></t<>	Methods	SRCC	PLCC	SRCC	PLCC	SRCC	PLCC	SRCC	PLCC	SRCC	PLCC	SRCC	PLCC
MinimalisticVQA(IX)	Weak Teachers												
FAST-VQA	MinimalisticVQA(VII)	0.861	0.859	0.740	0.784	0.843	0.841	0.757	0.813	0.775	0.779	0.817	0.830
DOVER	MinimalisticVQA(IX)	0.885	0.882	0.792	0.828	0.862	0.859	0.775	0.821	0.826	0.821	0.849	0.859
Q-Align 0.886 0.884 0.761 0.822 0.876 0.878 0.783 0.819 0.834 0.846 0.844 0.844 0.884 Weak-to-Strong Students Minimalistic VQA(IVI)-labeled 0.855 0.852 0.762 0.795 0.859 0.857 0.771 0.813 0.808 0.821 0.824 0.826 FAST-VQA-labeled 0.879 0.878 0.794 0.826 0.869 0.871 0.786 0.823 0.823 0.825 0.834 0.840 0.85 DOVER-labeled 0.877 0.869 0.780 0.813 0.870 0.875 0.792 0.829 0.819 0.831 0.849 0.85 Out of Distribution Datasets 0.876 0.794 0.824 0.873 0.874 0.874 0.874 0.797 0.828 0.830 0.831 0.831 0.88 0.89 0.831 0.88 0.89 0.89 0.89 0.89 0.88 0.88 0.830 0.831 0.834 0.89 <td>FAST-VQA</td> <td>0.880</td> <td>0.880</td> <td>0.781</td> <td>0.813</td> <td>0.859</td> <td>0.854</td> <td>0.826</td> <td>0.845</td> <td>0.730</td> <td>0.747</td> <td>0.838</td> <td>0.849</td>	FAST-VQA	0.880	0.880	0.781	0.813	0.859	0.854	0.826	0.845	0.730	0.747	0.838	0.849
Weak-to-Strong Students Minimalistic VQA(VII)-labeled 0.855 0.852 0.762 0.795 0.859 0.857 0.771 0.813 0.808 0.821 0.824 0.8 Minimalistic VQA(IX)-labeled 0.879 0.878 0.794 0.826 0.869 0.871 0.786 0.822 0.843 0.846 0.849 0.8 FAST-VQA-labeled 0.871 0.868 0.785 0.819 0.849 0.855 0.798 0.833 0.825 0.834 0.840 0.8 DOVER-labeled 0.878 0.876 0.794 0.824 0.873 0.880 0.813 0.875 0.792 0.829 0.819 0.831 0.84 0.8 Q-Align-labeled 0.881 0.876 0.794 0.824 0.873 0.880 0.781 0.825 0.833 0.853 0.831 0.88 0.879 0.834 0.874 0.871 0.825 0.833 0.851 0.8 0.8 0.8 0.8 0.8 0	DOVER	0.878	0.866	0.782	0.813	0.874	0.869	0.817	0.840	0.771	0.781	0.842	0.845
Minimalistic VQA(VII)-labeled 0.855 0.852 0.762 0.795 0.859 0.857 0.771 0.813 0.808 0.821 0.824 0.84 Minimalistic VQA(IX)-labeled 0.879 0.878 0.794 0.826 0.869 0.871 0.786 0.822 0.843 0.846 0.849 0.8 FAST-VQA-labeled 0.871 0.868 0.785 0.819 0.849 0.855 0.798 0.833 0.825 0.834 0.843 0.840 0.8 Q-Align-labeled 0.878 0.876 0.794 0.824 0.873 0.869 0.780 0.813 0.870 0.825 0.792 0.829 0.819 0.831 0.843 0.8 Q-Align-labeled 0.881 0.876 0.794 0.824 0.873 0.880 0.781 0.825 0.833 0.853 0.843 0.8 Supervised Student LIVE-YT-Gaming CGVDS LIVE-YT-HFR Waterlow-IVC-4K KVQ 0.821 0.841 0.82 </td <td>Q-Align</td> <td>0.886</td> <td>0.884</td> <td>0.761</td> <td>0.822</td> <td>0.876</td> <td>0.878</td> <td>0.783</td> <td>0.819</td> <td>0.834</td> <td>0.846</td> <td>0.844</td> <td>0.861</td>	Q-Align	0.886	0.884	0.761	0.822	0.876	0.878	0.783	0.819	0.834	0.846	0.844	0.861
Minimalistic VQA(IX)-labeled 0.879 0.878 0.794 0.826 0.869 0.871 0.786 0.822 0.843 0.846 0.849 0.85 FAST-VQA-labeled 0.871 0.868 0.785 0.819 0.849 0.855 0.798 0.833 0.825 0.834 0.840 0.8 DOVER-labeled 0.877 0.869 0.780 0.813 0.875 0.792 0.829 0.819 0.831 0.88 0.8 Q-Align-labeled 0.878 0.876 0.794 0.824 0.873 0.880 0.781 0.825 0.833 0.853 0.843 0.8 Supervised Student LIVE-YT-Gaming CGVDS LIVE-YT-HFR Waterloo-IVC-4K KVQ Overall # of videos SRCC PLCC SRCC PLCC SRCC PLCC SRCC PLC SRCC PLC SRCC PLC SRCC PLC SRCC PLC SRCC PLC SRC PLC SRCC PLC<	Weak-to-Strong Students												
FAST-VQA-labeled 0.871 0.868 0.785 0.819 0.849 0.855 0.798 0.833 0.825 0.834 0.840 0.850	MinimalisticVQA(VII)-labeled	0.855	0.852	0.762	0.795	0.859	0.857	0.771	0.813	0.808	0.821	0.824	0.833
DOVER-labeled 0.877 0.869 0.780 0.813 0.870 0.875 0.792 0.829 0.819 0.831 0.843 0.84 0.84 0.874 0.873 0.880 0.781 0.825 0.833 0.853 0.848 0.85 0.85	MinimalisticVQA(IX)-labeled	0.879	0.878	0.794	0.826	0.869	0.871	0.786	0.822	0.843	0.846	0.849	0.859
Q-Align-labeled 0.878 0.876 0.794 0.824 0.873 0.880 0.781 0.825 0.833 0.833 0.848 0.88 Supervised Student LSVQ-labeled 0.881 0.878 0.797 0.834 0.874 0.874 0.797 0.828 0.830 0.838 0.851 0.8 Out of Distribution Datasets LIVE-YT-Gaming CGVDS LIVE-YT-HFR Waterloo-IVC-4K KVQ Overall # of videos SRCC PLCC	FAST-VQA-labeled	0.871	0.868	0.785	0.819	0.849	0.855	0.798	0.833	0.825	0.834	0.840	0.850
Supervised Student LSVQ-labeled 0.881 0.878 0.797 0.834 0.874 0.874 0.797 0.828 0.830 0.838 0.851 0.8	DOVER-labeled	0.877	0.869	0.780	0.813	0.870	0.875	0.792	0.829	0.819	0.831	0.843	0.850
Supervised Student LSVQ-labeled 0.881 0.878 0.797 0.834 0.874 0.874 0.797 0.828 0.830 0.830 0.831 0.851 0.8 Out of Distribution Datasets LIVE-YT-Gaming CGVDS LIVE-YT-HFR Waterloo-IVC-4K KVQ Overall Methods SRCC PLCC	Q-Align-labeled	0.878	0.876	0.794	0.824	0.873	0.880	0.781	0.825	0.833	0.853	0.848	0.859
Hethods SRCC PLCC SRCC	Supervised Student												
# of videos		0.881	0.878	0.797	0.834	0.874	0.874	0.797	0.828	0.830	0.838	0.851	0.861
Methods SRCC PLCC SRCC	Out of Distribution Datasets	LIVE-Y	T-Gaming	CG	VDS	LIVE-Y	T-HFR	Waterlo	o-IVC-4K	K	VQ	Ove	rall
Weak Teachers Minimalistic VQA(VII) 0.596 0.682 0.681 0.733 0.061 0.130 0.275 0.338 0.604 0.659 0.490 0.55 Minimalistic VQA(IX) 0.686 0.746 0.797 0.816 0.301 0.388 0.459 0.502 0.615 0.661 0.574 0.6 FAST-VQA 0.631 0.677 0.725 0.747 0.326 0.415 0.327 0.363 0.518 0.526 0.486 0.5 DOVER 0.647 0.728 0.694 0.747 0.360 0.465 0.368 0.418 0.559 0.593 0.519 0.5 Q-Align 0.611 0.681 0.756 0.798 0.329 0.342 0.414 0.497 0.613 0.655 0.555 0.555 0.65 Weak-to-Strong Students Weak-to-Strong Students Weak-to-Strong Students Weak-to-Strong Students 0.612 0.638 0.717 0.718 0.773 0.318 0.386 0.412 0.6	# of videos		500	357		480		1,200		2,926		-	
Minimalistic VQA(VII) 0.596 0.682 0.681 0.733 0.061 0.130 0.275 0.338 0.604 0.659 0.490 0.55 Minimalistic VQA(IX) 0.686 0.746 0.797 0.816 0.301 0.388 0.459 0.502 0.615 0.661 0.574 0.6 FAST-VQA 0.631 0.677 0.725 0.747 0.326 0.415 0.327 0.363 0.518 0.526 0.486 0.5 DOVER 0.647 0.728 0.694 0.747 0.360 0.465 0.368 0.418 0.559 0.593 0.519 0.5 Q-Align 0.611 0.681 0.756 0.798 0.329 0.342 0.414 0.497 0.613 0.655 0.555 0.55 0.5 0.6 Weak-to-Strong Students Weak-to-Strong Students Winimalistic VQA(II)-labeled 0.632 0.717 0.718 0.773 0.318 0.386 0.3412 0.604 0.652 0.536 0.5<	Methods	SRCC	PLCC	SRCC	PLCC	SRCC	PLCC	SRCC	PLCC	SRCC	PLCC	SRCC	PLCC
Minimalistic VQA(IX) 0.686 0.746 0.797 0.816 0.301 0.388 0.459 0.502 0.615 0.661 0.574 0.6 FAST-VQA 0.631 0.677 0.725 0.747 0.326 0.415 0.327 0.363 0.518 0.526 0.486 0.5 DOVER 0.647 0.728 0.694 0.747 0.360 0.465 0.368 0.418 0.559 0.593 0.519 0.5 0.5 Q-Align 0.611 0.681 0.756 0.798 0.329 0.342 0.414 0.497 0.613 0.655 0.555 0.55 0.6 Weak-to-Strong Students Weak-to-Strong Students Weak-to-Strong Students 0.604 0.632 0.717 0.718 0.773 0.318 0.386 0.356 0.412 0.604 0.652 0.536 0.5 Minimalistic VQA(IX)-labeled 0.687 0.748 0.763 0.810 0.383 0.461 0.459 0.515 0.638 0.67	Weak Teachers												
Minimalistic VQA(IX) 0.686 0.746 0.797 0.816 0.301 0.388 0.459 0.502 0.615 0.661 0.574 0.6 FAST-VQA 0.631 0.677 0.725 0.747 0.326 0.415 0.327 0.363 0.518 0.526 0.486 0.5 DOVER 0.647 0.728 0.694 0.747 0.360 0.465 0.368 0.418 0.559 0.593 0.519 0.5 0.5 Q-Align 0.611 0.681 0.756 0.798 0.329 0.342 0.414 0.497 0.613 0.655 0.555 0.55 0.6 Weak-to-Strong Students Weak-to-Strong Students Weak-to-Strong Students 0.604 0.632 0.717 0.718 0.773 0.318 0.386 0.356 0.412 0.604 0.652 0.536 0.5 Minimalistic VQA(IX)-labeled 0.687 0.748 0.763 0.810 0.383 0.461 0.459 0.515 0.638 0.67	MinimalisticVOA(VII)	0.596	0.682	0.681	0.733	0.061	0.130	0.275	0.338	0.604	0.659	0.490	0.551
DOVER 0.647 0.728 0.694 0.747 0.360 0.465 0.368 0.418 0.559 0.593 0.519 0.55 Q-Align 0.611 0.681 0.756 0.798 0.329 0.342 0.414 0.497 0.613 0.655 0.555 0.655 0.555 0.6 Weak-to-Strong Students Winimalistic VQA(VII)-labeled 0.632 0.717 0.718 0.773 0.318 0.386 0.356 0.412 0.604 0.652 0.531 0.559 Minimalistic VQA(IX)-labeled 0.687 0.748 0.763 0.810 0.383 0.461 0.459 0.515 0.638 0.676 0.591 0.6 FAST-VQA-labeled 0.658 0.766 0.752 0.785 0.392 0.422 0.414 0.493 0.585 0.624 0.550 0.6 DOVER-labeled 0.662 0.758 0.752 0.809 0.449 0.482 0.435 0.519 0.574 0.627 0.554 0.6	2 . ,	0.686	0.746	0.797	0.816	0.301	0.388	0.459	0.502	0.615	0.661	0.574	0.622
Q-Align 0.611 0.681 0.756 0.798 0.329 0.342 0.414 0.497 0.613 0.655 0.555 0.6 Weak-to-Strong Students MinimalisticVQA(VII)-labeled 0.632 0.717 0.718 0.773 0.318 0.386 0.356 0.412 0.604 0.652 0.536 0.5 MinimalisticVQA(IX)-labeled 0.687 0.748 0.763 0.810 0.383 0.461 0.459 0.515 0.638 0.676 0.591 0.6 FAST-VQA-labeled 0.658 0.766 0.752 0.785 0.392 0.422 0.414 0.493 0.585 0.624 0.550 0.6 DOVER-labeled 0.662 0.758 0.752 0.809 0.449 0.482 0.435 0.519 0.574 0.627 0.554 0.6 Q-Align-labeled 0.671 0.738 0.744 0.785 0.437 0.480 0.450 0.525 0.620 0.668 0.581 0.6	FAST-VOA	0.631	0.677	0.725	0.747	0.326	0.415	0.327	0.363	0.518	0.526	0.486	0.512
Q-Align 0.611 0.681 0.756 0.798 0.329 0.342 0.414 0.497 0.613 0.655 0.555 0.6 Weak-to-Strong Students Minimalistic VQA(VII)-labeled 0.632 0.717 0.718 0.773 0.318 0.386 0.356 0.412 0.604 0.652 0.536 0.55 Minimalistic VQA(XIY)-labeled 0.687 0.748 0.763 0.810 0.383 0.461 0.459 0.515 0.638 0.676 0.591 0.6 FAST-VQA-labeled 0.658 0.766 0.752 0.785 0.392 0.422 0.412 0.493 0.585 0.624 0.550 0.6 DOVER-labeled 0.662 0.758 0.752 0.809 0.449 0.482 0.435 0.519 0.574 0.627 0.554 0.6 Q-Align-labeled 0.671 0.738 0.744 0.785 0.437 0.480 0.450 0.525 0.620 0.668 0.581 0.6	•	0.647			0.747	0.360			0.418		0.593		0.569
Weak-to-Strong Students Minimalistic VQA(VII)-labeled 0.632 0.717 0.718 0.773 0.318 0.386 0.356 0.412 0.604 0.652 0.536 0.5 Minimalistic VQA(IX)-labeled 0.687 0.748 0.763 0.810 0.383 0.461 0.459 0.515 0.638 0.676 0.591 0.6 FAST-VQA-labeled 0.658 0.766 0.752 0.785 0.392 0.422 0.414 0.493 0.585 0.624 0.550 0.6 DOVER-labeled 0.662 0.758 0.752 0.809 0.449 0.482 0.435 0.519 0.574 0.627 0.554 0.6 Q-Align-labeled 0.671 0.738 0.744 0.785 0.437 0.480 0.450 0.525 0.620 0.668 0.581 0.6	O-Align	0.611		0.756	0.798		0.342	0.414			0.655		0.606
Minimalistic VQA(VII)-labeled 0.632 0.717 0.718 0.773 0.318 0.386 0.356 0.412 0.604 0.652 0.536 0.5 Minimalistic VQA(IX)-labeled 0.687 0.748 0.763 0.810 0.383 0.461 0.459 0.515 0.638 0.676 0.591 0.6 FAST-VQA-labeled 0.658 0.766 0.752 0.785 0.392 0.422 0.414 0.493 0.585 0.624 0.550 0.6 DOVER-labeled 0.662 0.758 0.752 0.809 0.449 0.482 0.435 0.519 0.574 0.627 0.554 0.6 Q-Align-labeled 0.671 0.738 0.744 0.785 0.437 0.480 0.450 0.525 0.620 0.668 0.581 0.6	<u> </u>												
Minimalistic VQA(IX)-labeled 0.687 0.748 0.763 0.810 0.383 0.461 0.459 0.515 0.638 0.676 0.591 0.66 FAST-VQA-labeled 0.658 0.766 0.752 0.785 0.392 0.422 0.414 0.493 0.585 0.624 0.550 0.6 DOVER-labeled 0.662 0.758 0.752 0.809 0.449 0.482 0.435 0.519 0.574 0.627 0.554 0.6 Q-Align-labeled 0.671 0.738 0.744 0.785 0.437 0.480 0.450 0.525 0.620 0.688 0.581 0.6		0.632	0.717	0.718	0.773	0.318	0.386	0.356	0.412	0.604	0.652	0.536	0.593
FAST-VQA-labeled 0.658 0.766 <u>0.752</u> 0.785 0.392 0.422 0.414 0.493 0.585 0.624 0.550 0.6 DOVER-labeled 0.662 <u>0.758</u> <u>0.752</u> <u>0.809</u> <u>0.449</u> <u>0.482</u> 0.435 <u>0.519</u> 0.574 0.627 0.554 0.6 Q-Align-labeled <u>0.671</u> 0.738 0.744 0.785 0.437 0.480 0.450 0.525 <u>0.620</u> <u>0.668</u> <u>0.581</u> <u>0.68</u>													0.639
DOVER-labeled 0.662 0.758 0.752 0.809 0.449 0.482 0.435 0.519 0.574 0.627 0.554 0.6 Q-Align-labeled 0.671 0.738 0.744 0.785 0.437 0.480 0.450 0.525 0.620 0.668 0.581 0.6	• ' '												0.604
Q-Align-labeled <u>0.671</u> 0.738 0.744 0.785 0.437 0.480 0.450 0.525 <u>0.620</u> <u>0.668</u> <u>0.581</u> <u>0.6</u>	_												0.617
. <u></u>													0.636
Supervised Student	Supervised Student	- = -								- ===			

$$P(b^{(i)}, b^{(j)}) = 1 - P(b^{(j)}, b^{(i)}), \quad P(b^{(i)}, b^{(i)}) = 0.5.$$
 (30)

0.451

0.485

0.619

0.636

0.577

0.608

0.490

To evaluate a test video v_{test} , we compute its comparative probabilities against all anchor videos, forming the probability vector:

$$c = \left[P(b^{(1)}, v_{\text{test}}), P(b^{(2)}, v_{\text{test}}), \dots, P(b^{(5)}, v_{\text{test}}) \right].$$
 (31)

Next, we integrate this vector into the complete probability matrix:

0.643

0.713

0.713

0.770

0.451

$$M \in \mathbb{R}^{(5+1)\times(5+1)}, M = \begin{bmatrix} M_r & c \\ (1-c)^\top & 0.5 \end{bmatrix}.$$
 (32)

With this probability matrix, we estimate the final quality score using maximum a posteriori (MAP) (Tsukida et al., 2011) estimation under Thurstone's Case V model (Thurstone, 2017). This is formulated as the following convex optimization problem:

$$\begin{split} \arg \max_{\hat{q}} \sum_{i,j} M_{i,j} \log \left(\Phi(\hat{q}^{(i)} - \hat{q}^{(j)}) \right) \\ - \sum_{i} \frac{\hat{q}^{(i)}}{2}, \quad \text{s.t.} \sum_{i} \hat{q}^{(i)} = 0. \end{split} \tag{33}$$

Here, $\Phi(\cdot)$ denotes the standard normal cumulative distribution function, and the final score $\hat{q}^{(n+1)}$ corresponds to the estimated quality of the test video.

D More Details of Experimental Results

D.1 More Details of Weak-to-strong Generalization Effect

Table 4 presents the per-dataset results from the experiments described in Section 3.3 of the main paper. For in-domain benchmarks, the student model achieves performance comparable to its teachers, with slight improvements, demonstrating that our simple knowledge distillation approach effectively transfers quality assessment knowledge from weak to strong models. For OOD benchmarks, the student model shows substantial improvements over its teachers, highlighting a pronounced weak-to-strong generalization effect.

REFERENCES

Joao Carreira and Andrew Zisserman. Quo vadis, action recognition? a new model and the kinetics dataset. In *proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 6299–6308, 2017.

Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *IEEE Conference on Computer Vision and Pattern Recognition*, pp. 248–255, 2009.

Christoph Feichtenhofer, Haoqi Fan, Jitendra Malik, and Kaiming He. Slowfast networks for video recognition. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 6202–6211, 2019.

David Hasler and Sabine E Suesstrunk. Measuring colorfulness in natural images. In *Human vision and electronic imaging VIII*, volume 5007, pp. 87–95. SPIE, 2003.

Vlad Hosu, Franz Hahn, Mohsen Jenadeleh, Hanhe Lin, Hui Men, Tamás Szirányi, Shujun Li, and Dietmar Saupe. The konstanz natural video database (konvid-1k). In 2017 Ninth international Conference on Quality of Multimedia experience, pp. 1–6, 2017.

ITU-T P.910. Subjective video quality assessment methods for multimedia applications, 2008. URL https://www.itu.int/rec/dologin_pub.asp?lang=e&id=T-REC-P.910-200804-S!!PDF-E&type=items.

Zhuoran Li, Zhengfang Duanmu, Wentao Liu, and Zhou Wang. Avc, hevc, vp9, avs2 or av1?—a comparative study of state-of-the-art video encoders on 4k videos. In *Image Analysis and Recognition: 16th International Conference, ICIAR 2019, Waterloo, ON, Canada, August 27–29, 2019, Proceedings, Part I 16*, pp. 162–173. Springer, 2019.

Ze Liu, Jia Ning, Yue Cao, Yixuan Wei, Zheng Zhang, Stephen Lin, and Han Hu. Video swin transformer. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 3202–3211, 2022.

Yiting Lu, Xin Li, Yajing Pei, Kun Yuan, Qizhi Xie, Yunpeng Qu, Ming Sun, Chao Zhou, and Zhibo Chen. Kvq: Kwai video quality assessment for short-form videos. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 25963–25973, 2024.

Pavan C Madhusudana, Xiangxu Yu, Neil Birkbeck, Yilin Wang, Balu Adsumilli, and Alan C Bovik. Subjective and objective quality assessment of high frame rate videos. *IEEE Access*, 9:108069–108082, 2021.

Niranjan D Narvekar and Lina J Karam. A no-reference image blur metric based on the cumulative probability of blur detection (cpbd). *IEEE Transactions on Image Processing*, 20(9):2678–2683, 2011.

- Juergen Pandel. Measuring of flickering artifacts in predictive coded video sequences. In 2008 Ninth International Workshop on Image Analysis for Multimedia Interactive Services, pp. 231–234. IEEE, 2008.
- Eli Peli. Contrast in complex images. *JOSA A*, 7(10):2032–2040, 1990.
 - Piotr Romaniak, Lucjan Janowski, Mikolaj Leszczuk, and Zdzislaw Papir. Perceptual quality assessment for h. 264/avc compression. In 2012 IEEE Consumer Communications and Networking Conference, pp. 597–602. IEEE, 2012.
 - Avinab Saha, Yu-Chih Chen, Chase Davis, Bo Qiu, Xiaoming Wang, Rahul Gowda, Ioannis Katsavounidis, and Alan C Bovik. Study of subjective and objective quality assessment of mobile cloud gaming videos. *IEEE Transactions on Image Processing*, 32:3295–3310, 2023.
 - Zeina Sinno and Alan Conrad Bovik. Large-scale study of perceptual video quality. *IEEE Transactions on Image Processing*, 28(2):612–627, 2018.
 - Wei Sun, Wen Wen, Xiongkuo Min, Long Lan, Guangtao Zhai, and Kede Ma. Analysis of video quality datasets via design of minimalistic video quality models. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2024.
 - Louis L Thurstone. A law of comparative judgment. In Scaling, pp. 81–92. Routledge, 2017.
 - Kristi Tsukida, Maya R Gupta, et al. How to analyze paired comparison data. *Department of Electrical Engineering University of Washington, Tech. Rep. UWEETR-2011-0004*, 1, 2011.
 - Vassilios Vonikakis, Ramanathan Subramanian, Jonas Arnfred, and Stefan Winkler. A probabilistic approach to people-centric photo selection and sequencing. *IEEE Transactions on Multimedia*, 19(11):2609–2624, 2017.
 - Yilin Wang, Sasi Inguva, and Balu Adsumilli. Youtube ugc dataset for video compression research. In 2019 IEEE 21st International Workshop on Multimedia Signal Processing, pp. 1–5. IEEE, 2019
 - Haoning Wu, Chaofeng Chen, Jingwen Hou, Liang Liao, Annan Wang, Wenxiu Sun, Qiong Yan, and Weisi Lin. Fast-vqa: Efficient end-to-end video quality assessment with fragment sampling. In *European Conference on Computer Vision*, pp. 538–554. Springer, 2022.
 - Haoning Wu, Erli Zhang, Liang Liao, Chaofeng Chen, Jingwen Hou, Annan Wang, Wenxiu Sun, Qiong Yan, and Weisi Lin. Exploring video quality assessment on user generated contents from aesthetic and technical perspectives. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 20144–20154, 2023a.
 - Haoning Wu, Zicheng Zhang, Weixia Zhang, Chaofeng Chen, Liang Liao, Chunyi Li, Yixuan Gao, Annan Wang, Erli Zhang, Wenxiu Sun, et al. Q-align: Teaching lmms for visual scoring via discrete text-defined levels. *arXiv preprint arXiv:2312.17090*, 2023b.
 - Zhenqiang Ying, Maniratnam Mandal, Deepti Ghadiyaram, and Alan Bovik. Patch-vq:'patching up'the video quality problem. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 14019–14029, 2021.
 - Xiangxu Yu, Zhengzhong Tu, Zhenqiang Ying, Alan C Bovik, Neil Birkbeck, Yilin Wang, and Balu Adsumilli. Subjective quality assessment of user-generated content gaming videos. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pp. 74–83, 2022.